



PATENT ABSTRACTS OF JAPAN

(11) Publication number: **2002182682 A**(43) Date of publication of application: **26.06.02**

(51) Int. Cl. **G10L 15/06**
G10L 15/14
G10L 21/04
G10L 13/06
G10L 15/02

(21) Application number: **2000382371**(22) Date of filing: **15.12.00**(71) Applicant: **SHARP CORP**(72) Inventor: **YAMAGUCHI KOICHI**
HACHIMAN YOICHIRO

(54) **SPEAKER CHARACTERISTIC EXTRACTOR, SPEAKER CHARACTERISTIC EXTRACTION METHOD, SPEECH RECOGNIZER, SPEECH SYNTHESIZER AS WELL AS PROGRAM RECORDING MEDIUM**

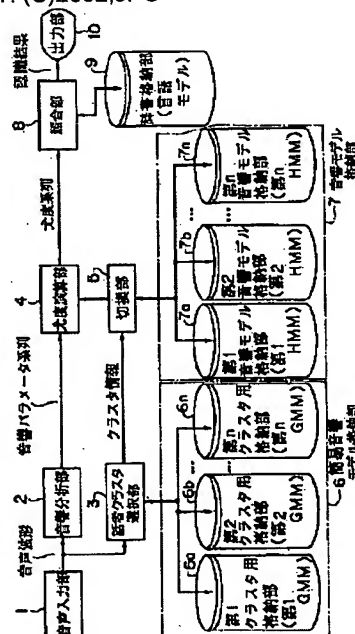
extracted with good accuracy by taking the speakers' habits into consideration from a smaller quantity of the utterance data as the distances between the respective learning speakers.

COPYRIGHT: (C)2002,JPO

(57) Abstract:

PROBLEM TO BE SOLVED: To extract speaker characteristics with good accuracy from a smaller quantity of utterance data.

SOLUTION: Acoustic models are stored in a first acoustic model storage section 7a to an n-th acoustic model storage section 7n by each of n pieces of speaker clusters in the acoustic model storage sections 7. The vocal tract length normalization coefficient α determined by estimating likelihood by equation (a) according to a reference of maximizing the likelihood of the acoustic models of learning speakers for the acoustic models of all the learning speakers by using a nonlinear frequency warping obtained by applying a correction factor β to vocal tract length normalization coefficient α is used for clustering of the learning speakers of this case as the distance between the respective learning speakers. The distances between the respective learning speakers are set in accordance with the information on the vocal tract lengths which are the fluctuating factors of the physiological characteristics and the correction information of the ways and habits of the utterance, by which the learning speakers are clustered with the speaker characteristics



【特許請求の範囲】

【請求項 1】 入力音声から、標準話者の音声のスペクトルに対して上記入力音声のスペクトルの周波数軸を伸縮する際の伸縮係数 α を話者特徴として抽出する話者特徴抽出装置において、

音声のスペクトルにおける広母音の第 2 フォルマンの存在領域以下の低い周波数領域で上記伸縮係数 α に補正係数 β を乗じて上記伸縮係数 α に対して部分的に補正を行った非線形周波数ワーピング関数を用いて、標準話者の音声パターンに対して入力話者の音声パターンの尤度を最大にするという基準に従って、上記伸縮係数 α を求める伸縮係数取得手段を備えたことを特徴とする話者特徴抽出装置。

【請求項 2】 請求項 1 に記載の話者特徴抽出装置において、

上記伸縮係数取得手段は、上記非線形周波数ワーピング関数を用いて、標準話者の音声パターンに対して入力話者の音声パターンの尤度を最大にするという基準に従って、上記補正係数 β をも求めるようになっていることを特徴とする話者特徴抽出装置。

【請求項 3】 標準話者の音声スペクトルに対して学習話者の音声スペクトルの周波数軸を伸縮する際の伸縮係数 α に基づいて上記学習話者をクラスタリングし、各話者クラスに属する学習話者群の音声パターンに基づいて作成された音響モデルを、各話者クラス別の音響モデル格納部に格納した音声認識装置であって、上記伸縮係数 α は、音声のスペクトルにおける広母音の第 2 フォルマンの存在領域以下の低い周波数領域で上記伸縮係数 α に補正係数 β を乗じて上記伸縮係数 α に対して部分的に補正を行った非線形周波数ワーピング関数を用いて、上記標準話者の音声パターンに対して学習話者の音声パターンの尤度を最大にするという基準に従って求められていることを特徴とする音声認識装置。

【請求項 4】 請求項 3 に記載の音声認識装置において、

上記学習話者のクラスタリングは、上記伸縮係数 α と補正係数 β との 2 次元平面に対して行われており、上記補正係数 β は、上記非線形周波数ワーピング関数を用いて、上記標準話者の音声パターンに対して学習話者の音声パターンの尤度を最大にするという基準に従って求められていることを特徴とする音声認識装置。

【請求項 5】 請求項 3 あるいは請求項 4 に記載の音声認識装置において、

上記話者クラスは、

上記伸縮係数 α による 1 次元空間、あるいは、上記伸縮係数 α と補正係数 β とによる 2 次元空間を、所定のクラス数にクラスタリングした初期クラスと、

上記各初期クラスの境界を含んで上記各初期クラスにオーバーラップするオーバーラップクラスで構成されていることを特徴とする音声認識装置。

【請求項 6】 標準話者の音声スペクトルに対して入力話者の音声スペクトルの周波数軸を伸縮する際の伸縮係数 α を用いて入力話者の音声スペクトルの周波数軸を伸縮することによって上記入力話者の音声を正規化する正規化手段を有する音声認識装置において、

上記正規化手段は、

音声のスペクトルにおける広母音の第 2 フォルマンの存在領域以下の低い周波数領域で上記伸縮係数 α に補正係数 β を乗じて上記伸縮係数 α に対して部分的に補正を行った非線形周波数ワーピング関数を用いて、標準話者の音声パターンに対して入力話者の音声パターンの尤度を最大にするという基準に従って、上記伸縮係数 α と補正係数 β とを推定する周波数ワーピング関数推定手段

と、

上記推定された伸縮係数 α と補正係数 β を係数とする上記非線形周波数ワーピング関数を用いて、上記入力話者の音声スペクトルの周波数軸を伸縮する周波数ワープ手段で構成されていることを特徴とする音声認識装置。

【請求項 7】 標準話者の音声スペクトルに対して入力話者の音声スペクトルの周波数軸を伸縮する際の伸縮係数 α を用いて音声のスペクトルの周波数軸を伸縮することによって音響モデルを入力話者に話者適応させる話者適応手段を有する音声認識装置において、

上記話者適応手段は、

音声のスペクトルにおける広母音の第 2 フォルマンの存在領域以下の低い周波数領域で上記伸縮係数 α に補正係数 β を乗じて上記伸縮係数 α に対して部分的に補正を行った非線形周波数ワーピング関数を用いて、標準話者の音声パターンに対して入力話者の音声パターンの尤度を最大にするという基準に従って、上記伸縮係数 α と補正係数 β とを推定する周波数ワーピング関数推定手段

と、

上記推定された伸縮係数 α の逆数と補正係数 β の逆数とを係数とする上記非線形周波数ワーピング関数を用いて、上記音響モデルの周波数軸を伸縮する周波数ワープ手段で構成されていることを特徴とする音声認識装置。

【請求項 8】 請求項 6 あるいは請求項 7 に記載の音声認識装置において、

上記周波数ワーピング関数推定手段は、上記入力話者の音声パターンの代わりに、標準話者の音響モデルを上記入力話者の音声パターンに話者適応させて作成された入力話者用の適応音響モデルを用いるようになっていることを特徴とする音声認識装置。

【請求項 9】 請求項 4 乃至請求項 8 の何れか一つに記載の音声認識装置において、

上記補正係数 β は、音響モデルの状態や音素等のサブワード単位に求められ、上記サブワード毎に決定されていることを特徴とする音声認識装置。

【請求項 10】 入力話者の音声スペクトルに対して標準話者の音声スペクトルの周波数軸を伸縮する際の伸縮

係数 α を用いて音声のスペクトルの周波数軸を伸縮することによって、標準話者の音声素片を接続して成る合成音声の声質を発話者の声質に変換する声質変換手段を有する音声合成装置において、

上記声質変換手段は、

音声のスペクトルにおける広母音の第2フォルマントの存在領域以下の低い周波数領域で上記伸縮係数 α に補正係数 β を乗じて上記伸縮係数 α に対して部分的に補正を行った非線形周波数ワーピング関数を用いて、標準話者の音声パターンに対して上記発話者の音声パターンの尤度を最大にするという基準に従って、上記伸縮係数 α と補正係数 β とを推定する周波数ワーピング関数推定手段と、

上記推定された伸縮係数 α の逆数と補正係数 β の逆数とを係数とする上記非線形周波数ワーピング関数を用いて、上記音声素片の周波数軸を伸縮する周波数ワーピング手段で構成されていることを特徴とする音声合成装置。

【請求項11】 請求項10に記載の音声合成装置において、

上記周波数ワーピング関数推定手段は、上記補正係数 β を、音響モデルの状態や音素等のサブワード単位で求め、そのサブワード毎に推定するようになっていることを特徴とする音声合成装置。

【請求項12】 入力音声から、標準話者の音声のスペクトルに対して上記入力音声のスペクトルの周波数軸を伸縮する際の伸縮係数 α を話者特徴として抽出する話者特徴抽出方法において、

音声のスペクトルにおける広母音の第2フォルマントの存在領域以下の低い周波数領域で上記伸縮係数 α に補正係数 β を乗じて上記伸縮係数 α に対して部分的に補正を行った非線形周波数ワーピング関数を用いて、標準話者の音声パターンに対して入力話者の音声パターンの尤度を最大にするという基準に従って、上記伸縮係数 α を求めることを特徴とする話者特徴抽出方法。

【請求項13】 コンピュータを、

請求項1における上記伸縮係数取得手段として機能させる話者特徴抽出処理プログラムが記録されたことを特徴とするコンピュータ読出し可能なプログラム記録媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】この発明は、標準話者の音声スペクトルに対する入力音声スペクトルの周波数軸の線形伸縮係数を話者特徴として抽出する話者特徴抽出装置および話者特徴抽出方法、その抽出方法を用いた音声認識装置、音声合成装置、並びに、話者特徴抽出処理プログラムを記録したプログラム記録媒体に関する。

【0002】

【従来の技術】従来より、隠れマルコフモデル(Hidden Markov Model: 以下、HMMと言う)を用いた音声認識方法の開発が近年盛んに行われている。このHMMは、

大量の音声データから得られる音声の統計的特徴を確率的にモデル化したものであり、このHMMを用いた音声認識方法の詳細は、中川聖一著「確率モデルによる音声認識」(電子情報通信学会)に詳しい。このHMMに基づく話者適応や話者正規化に関する研究が行われており、特に声道長に基づく話者正規化が盛んに研究されて効果が出ている。声道長の差は話者間の主な変動要因であり、声道長は従来の話者適応法に比べて1個のパラメータあるいは極めて少ないパラメータで音声の特徴を表現できることから、声道長にはより少量の学習データで効率良く正規化できるというメリットがある。

【0003】ところで、標準話者の音声パターンに対する入力話者の音声サンプルの尤度を最大にするという基準(最尤推定)に従って、上記音声サンプルにおける周波数軸の線形伸縮係数 α (声道長正規化係数)を求める。そして、この声道長正規化係数 α を用いて入力話者の音声サンプルの周波数軸を線形伸縮することで話者正規化する技術が提案されている(例えば、AT&T Bell Labs. Li Lee, Richard C. Rose, "Speaker Normalization using Efficient Frequency Warping Procedures", pp. 353-356 ICASSP96 (1996))。また、特開平11-327592号公報においては、声道を前室と後室との2つの室に分け、各室に対応した2つの周波数軸線形伸縮係数 α を用いて話者正規化する技術が開示されている。

【0004】尚、上記話者適応は標準となる音響モデルを入力話者に対して適応(つまり正規化)させる技術であり、話者正規化とは表裏一体の関係にある。

【0005】また、話者クラスタリングを用いた音声認識方法がある。この音声認識方法においては、学習話者間の距離を定義して学習話者をクラスタリングしておき、クラスタ毎にそのクラスタに属する学習話者群の音声データを用いて音響モデルを作成する。そして、認識時には、入力音声に最適なクラスタを選択し、そのクラスタの音響モデルを用いて認識処理を行うのである。その場合における学習話者間の距離として上記声道長の周波数軸線形伸縮係数を用いる音声認識装置が提案されている(特開平11-175090号公報)。この公報においては、声道を前室と後室との2つの室に分け、各室に対応した2つの周波数軸線形伸縮係数を用いて学習話者をクラスタリングするようにしている。

【0006】さらに、音声合成における声質変換に関する従来技術として、音声認識の話者適応技術を用いてスペクトルの写像を行なう方法が提案されている。例えば、ベクトル量子化(VQ)コードブックマッピング法をベースとした話者適応技術を用いる方法(特開平7-104792号公報)や、VFS(Vector Field Smoothing)法をベースとした話者適応技術を用いる方法(橋本誠、樋口宣男: "話者選択と移動ベクトル場平滑化を用いた声質変換のためのスペクトル写像", 信学技報, SP95-1, p. p. 1-8, May 1995)等がある。

【0007】

【発明が解決しようとする課題】しかしながら、上記従来の声道長に基づく話者適応や話者正規化には、以下のような問題がある。すなわち、声道長に基づく話者適応や話者正規化は極めて少ないパラメータ数で音声の特徴を表現できるとは言うものの、話者内変動と呼ばれるその時の発話の仕方や癖等の影響を受け易い。したがって、必ずしも少ない学習サンプルから安定して声道長を抽出できるとは限らない。そして、現在一般に用いられている音声データベースには、身長、出身地、性別、年齢等の情報しか付与されていないのである。

【0008】実際の声道長はMRI(磁気共鳴画像診断装置)で測定しなければ分からないため、現時点においては直ちに真の声道長を知るのは困難な状況にある。上記特開平11-327592号公報および特開平11-175090号公報では、声道パラメータを得るために入力音声のフォルマント周波数を用いている。しかしながら、一般的にフォルマント周波数を全自動で求めることは困難であり、上記特開平11-327592号公報に開示された線形伸縮係数を用いた話者正規化方法や上記特開平11-175090号公報に開示された線形伸縮係数を用いた音声認識装置では、実時間性に欠けるという問題がある。

【0009】さらに、発声の仕方や一部の発音器官の形状による違いもフォルマント周波数のずれとなって現れる。そのため、少ない発声データから声道長を正規化するのは一般には困難である。また、声道長の伸縮(すなわち周波数軸のワーピング)を線形関数やそれに類似した関数で表現している。そのために、全区間に対して一様に周波数ワープを作用させることになり、声道長の差の影響を受け難い音素や無音部まで正規化(すなわち変形)されてしまうという問題もある。

【0010】また、上記特開平11-175090号公報のごとく、話者クラスタリングを用いた音声認識のアプローチも盛んに試みられているが、大きな性能改善は達成できていない。不特定話者(SI)音響モデル(すなわち男女共通の音響モデル)をベースラインとすると、男女別(GD)音響モデルは最もシンプルながら性能向上量が最も大きい。しかしながら、話者クラスタによって更なる細分化(クラスタ化)を行っても効果は薄いという報告がなされており、その場合における単語誤り率(WER: Word Error Rate)の削減は10%~20%程度に留まっている。これは、話者間の距離を定義する適当な尺度がないために上手くクラスタリングできなかったり、クラスタを増やすと1つのクラスタ当りの学習話者数が少なくなってロバスト性に欠けたりするためである。

【0011】さらに、何れの音響モデルの場合も、各話者クラスタの境界領域では学習サンプルが希薄だったり段差ができていたりしているため上手く学習されていない。

したがって、入力話者が各クラスタの境界付近に位置する場合には、認識率が劣化するという問題(所謂、hard decision問題)が生ずることになる。尚、個々の学習話者の音響モデル間の距離でクラスタリングを行った場合は、クラスタを木構造にし、入力話者が二つのクラスタの境界付近に位置する場合は上記2つのクラスタの上位ノードのクラスタの音響モデルを採用する方法もある。しかしながら、この方法の場合には、二つのクラスタの境界付近に位置する入力話者に対しては上位ノードの音響モデルを使用するためによりブロードな音響モデルとなってしまう、高い認識率は得にくいのである。

【0012】以上のごとく、上記話者適応(話者正規化)においては少ない発声データから音響モデルを精度良く適応できないため、誤り率を半減させるためには数十単語以上の発声データが必要となり、学習話者に負担を強いることになるという問題がある。また、音声合成における声質変換の場合にも、同様に少ない発声データからは精度良く声質が得られないという問題がある。

【0013】そこで、この発明の目的は、より少ない発声データから精度良く話者特徴を抽出できる話者特徴抽出装置および話者特徴抽出方法、その抽出方法を用いた音声認識装置、音声合成装置、並びに、話者特徴抽出処理プログラムを記録したプログラム記録媒体を提供することにある。

【0014】

【課題を解決するための手段】上記目的を達成するため、第1の発明は、入力音声から、標準話者の音声のスペクトルに対して上記入力音声のスペクトルの周波数軸を伸縮する際の伸縮係数 α を話者特徴として抽出する話者特徴抽出装置において、音声のスペクトルにおける広母音の第2フォルマントの存在領域以下の低い周波数領域で上記伸縮係数 α に補正係数 β を乗じて上記伸縮係数 α に対して部分的に補正を行った非線形周波数ワーピング関数を用いて、標準話者の音声パターンに対して入力話者の音声パターンの尤度を最大にするという基準に従って、上記伸縮係数 α を求める伸縮係数取得手段を備えたことを特徴としている。

【0015】上記構成によれば、音声のスペクトルにおける広母音の第2フォルマントの存在領域以下の低い周波数領域で伸縮係数 α に補正係数 β を乗じて補正を行った非線形周波数ワーピング関数を用いて、最尤推定によって上記伸縮係数 α を求め、その求められた伸縮係数 α をもって話者特徴としている。したがって、生理的な特徴の変動要因である声道長の情報に対して発声の仕方や癖による影響の補正が行われて、より話者に適合した特徴が抽出される。その際に、話者の発声の仕方や癖による影響が予め補正されている。したがって、上記発声の仕方や癖を表す発声データを必要とはせず、より少量の発声データから精度良く話者特徴が抽出される。

【0016】また、上記第1の発明の話者特徴抽出装置

は、上記伸縮係数取得手段を、上記非線形周波数ワーピング関数を用いて、標準話者の音声パターンに対して入力話者の音声パターンの尤度を最大にするという基準に従って、上記補正係数 β をも求めるように成すことが望ましい。

【0017】上記構成によれば、話者特徴として、上記伸縮係数 α に加えて、広母音の第2フォルマントの存在領域以下の低い周波数領域に対する補正係数 β もが抽出される。したがって、さらに話者に適合した特徴が抽出される。

【0018】また、第2の発明は、標準話者の音声スペクトルに対して学習話者の音声スペクトルの周波数軸を伸縮する際の伸縮係数 α に基づいて上記学習話者をクラスタリングし、各話者クラスタに属する学習話者群の音声パターンに基づいて作成された音響モデルを各話者クラスタ別の音響モデル格納部に格納した音声認識装置であって、上記伸縮係数 α は、音声のスペクトルにおける広母音の第2フォルマントの存在領域以下の低い周波数領域で上記伸縮係数 α に補正係数 β を乗じて上記伸縮係数 α に対して部分的に補正を行った非線形周波数ワーピング関数を用いて、上記標準話者の音声パターンに対して学習話者の音声パターンの尤度を最大にするという基準に従って求められていることを特徴としている。

【0019】上記構成によれば、学習話者をクラスタリングする際における各学習話者間の距離として、上記非線形周波数ワーピング関数を用いて最尤推定によって求めた上記伸縮係数 α が用いられる。こうして、生理的な特徴の変動要因である声道長の情報に対して発声の仕方や癖による影響の補正が行われて、より学習話者の音声パターンに適合した距離を用いて上記クラスタリングが行われる。その際に、話者の発声の仕方や癖による影響が予め補正されているので、上記発声の仕方や癖を表す発声データを必要とせず、より少量の発声データから学習話者間の距離が得られるのである。

【0020】また、上記第2の発明の音声認識装置は、上記学習話者のクラスタリングを、上記伸縮係数 α と補正係数 β との2次元平面に対して行うようになっており、上記補正係数 β を、上記非線形周波数ワーピング関数を用いて、上記標準話者の音声パターンに対して学習話者の音声パターンの尤度を最大にするという基準に従って求めることが望ましい。

【0021】上記構成によれば、上記学習話者間の距離として、上記伸縮係数 α に加えて、上記補正係数 β も用いられる。したがって、さらに話者に適合した距離を用いてクラスタリングが行われる。

【0022】また、上記第2の発明の音声認識装置は、上記話者クラスタを、上記伸縮係数 α による1次元空間あるいは上記伸縮係数 α と補正係数 β とによる2次元空間を、所定のクラスタ数にクラスタリングした初期クラスタと、上記各初期クラスタの境界を含んで上記各初期

クラスタにオーバーラップするオーバーラップクラスタとで構成することが望ましい。

【0023】上記構成によれば、話者クラスタは、所定数の初期クラスタと上記各初期クラスタにオーバーラップするオーバーラップクラスタとで構成されている。したがって、学習サンプルが希薄だったり段差がき易い上記各初期クラスタの境界領域は、何れかのオーバーラップクラスタに含まれることになり、上記各初期クラスタの境界領域において認識率が劣化するという「hard decision問題」が解消される。

【0024】また、第3の発明は、標準話者の音声スペクトルに対して入力話者の音声スペクトルの周波数軸を伸縮する際の伸縮係数 α を用いて入力話者の音声スペクトルの周波数軸を伸縮することによって上記入力話者の音声正規化する正規化手段を有する音声認識装置において、上記正規化手段は、音声のスペクトルにおける広母音の第2フォルマントの存在領域以下の低い周波数領域で上記伸縮係数 α に補正係数 β を乗じて上記伸縮係数 α に対して部分的に補正を行った非線形周波数ワーピング関数を用いて、標準話者の音声パターンに対して入力話者の音声パターンの尤度を最大にするという基準に従って、上記伸縮係数 α と補正係数 β とを推定する周波数ワーピング関数推定手段と、上記推定された伸縮係数 α と補正係数 β とを係数とする上記非線形周波数ワーピング関数を用いて、上記入力話者の音声スペクトルの周波数軸を伸縮する周波数ワープ手段で構成されていることを特徴としている。

【0025】上記構成によれば、周波数ワーピング関数推定手段によって、上記非線形周波数ワーピング関数を用いて、上記伸縮係数 α と補正係数 β とが最尤推定される。そして、周波数ワープ手段によって、推定された α と β とを係数とする上記非線形周波数ワーピング関数を用いて入力話者が正規化される。こうして、生理的な特徴の変動要因である声道長の情報に対して発声の仕方や癖による影響の補正が行われて、より標準話者の音声スペクトルに近づくように話者の正規化が行われる。その際に、話者の発声の仕方や癖による影響が予め補正されているので、上記発声の仕方や癖を表す発声データを必要とせず、より少量の発声データに基づいて話者正規化が行われるのである。

【0026】また、第4の発明は、標準話者の音声スペクトルに対して入力話者の音声スペクトルの周波数軸を伸縮する際の伸縮係数 α を用いて音声のスペクトルの周波数軸を伸縮することによって音響モデルを入力話者に話者適応させる話者適応手段を有する音声認識装置において、上記話者適応手段は、音声のスペクトルにおける広母音の第2フォルマントの存在領域以下の低い周波数領域で上記伸縮係数 α に補正係数 β を乗じて上記伸縮係数 α に対して部分的に補正を行った非線形周波数ワーピング関数を用いて、標準話者の音声パターンに対して入

力話者の音声パターンの尤度を最大にするという基準に従って、上記伸縮係数 α と補正係数 β とを推定する周波数ワーピング関数推定手段と、上記推定された伸縮係数 α の逆数と補正係数 β の逆数とを係数とする上記非線形周波数ワーピング関数を用いて、上記音響モデルの周波数軸を伸縮する周波数ワープ手段で構成されていることを特徴としている。

【0027】上記構成によれば、周波数ワーピング関数推定手段によって、上記非線形周波数ワーピング関数を用いて、上記伸縮係数 α と補正係数 β とが最尤推定される。そして、周波数ワープ手段によって、この推定された伸縮係数 α の逆数と補正係数 β の逆数とを係数とする上記非線形周波数ワーピング関数を用いて、音響モデルが入力話者に話者適応される。こうして、生理的な特徴の変動要因である声道長の情報に対して発声の仕方や癖による影響の補正が行われて、より入力話者の音声スペクトルに近づくように話者適応が行われる。その際に、話者の発声の仕方や癖による影響が予め補正されているので、上記発声の仕方や癖を表す発声データを必要とせず、より少量の発声データに基づいて話者適応が行われるのである。

【0028】また、上記第3の発明あるいは第4の発明の音声認識装置は、上記周波数ワーピング関数推定手段を、上記入力話者の音声パターンの代わりに、標準話者の音響モデルを上記入力話者の音声パターンに話者適応させて作成された入力話者用の適応音響モデルを用いるように成すことが望ましい。

【0029】上記構成によれば、上記伸縮係数 α と補正係数 β との推定に際して、入力話者の音声パターンそのものではなく、標準話者の音響モデルを入力話者の音声パターンに話者適応させた適応音響モデルを用いるので、入力話者の音声パターン数が少ない場合でも対処可能となる。さらに、上記適応音響モデルの状態毎に補正係数 β を制御して、話者の発声の仕方や癖による入力音声パターンのずれを木目細かく補正することが可能になる。

【0030】また、上記第2の発明乃至第4の発明の何れか一つの発明の音声認識装置は、上記補正係数 β を、音響モデルの状態や音素等のサブワード単位に求め、上記サブワード毎に決定することが望ましい。

【0031】上記構成によれば、上記補正係数 β がサブワード単位に変更されて、話者の発声の仕方や癖による入力音声パターンのずれが木目細かく補正される。

【0032】また、第5の発明は、標準話者の音声スペクトルに対して入力話者の音声スペクトルの周波数軸を伸縮する際の伸縮係数 α を用いて音声のスペクトルの周波数軸を伸縮することによって、標準話者の音声素片が接続されて成る合成音声の声質を発話者の声質に変換する声質変換手段を有する音声合成装置において、上記声質変換手段は、音声のスペクトルにおける広母音の第2

フォルマントの存在領域以下の低い周波数領域で上記伸縮係数 α に補正係数 β を乗じて上記伸縮係数 α に対して部分的に補正を行った非線形周波数ワーピング関数を用いて、標準話者の音声パターンに対して上記発話者の音声パターンの尤度を最大にするという基準に従って、上記伸縮係数 α と補正係数 β とを推定する周波数ワーピング関数推定手段と、上記推定された伸縮係数 α の逆数と補正係数 β の逆数とを係数とする上記非線形周波数ワーピング関数を用いて、上記音声素片の周波数軸を伸縮する周波数ワープ手段で構成されていることを特徴としている。

【0033】上記構成によれば、周波数ワーピング関数推定手段によって、上記非線形周波数ワーピング関数を用いて、最尤推定によって上記伸縮係数 α と補正係数 β とが推定される。そして、周波数ワープ手段によって、推定された伸縮係数 α の逆数と補正係数 β の逆数とを係数とする上記非線形周波数ワーピング関数を用いて、標準話者の音声素片の周波数軸が伸縮される。こうして、生理的な特徴の変動要因である声道長の情報に対して発声の仕方や癖による影響の補正が行われて、上記合成音声の声質が発話者の声質に変換される。その際に、話者の発声の仕方や癖による影響が予め補正されているので、上記発声の仕方や癖を表す発声データを必要とせず、より少量の発声データに基づいて声質変換が行われる。

【0034】また、上記第5の発明の音声合成装置は、上記周波数ワーピング関数推定手段を、音響モデルの状態や音素等のサブワード単位に上記補正係数 β を求め、上記サブワード毎に上記補正係数 β を推定するように成すことが望ましい。

【0035】上記構成によれば、上記補正係数 β がサブワード単位に変更されて、発話者の発声の仕方や癖による入力音声パターンのずれが木目細かく補正される。

【0036】また、第6の発明は、入力音声から、標準話者の音声のスペクトルに対して上記入力音声のスペクトルの周波数軸を伸縮する際の伸縮係数 α を話者特徴として抽出する話者特徴抽出方法において、音声のスペクトルにおける広母音の第2フォルマントの存在領域以下の低い周波数領域で上記伸縮係数 α に補正係数 β を乗じて上記伸縮係数 α に対して部分的に補正を行った非線形周波数ワーピング関数を用いて、標準話者の音声パターンに対して入力話者の音声パターンの尤度を最大にするという基準に従って、上記伸縮係数 α を求めることを特徴としている。

【0037】上記構成によれば、上記非線形周波数ワーピング関数を用いて、最尤推定によって伸縮係数 α が話者特徴として求められる。したがって、生理的な特徴の変動要因である声道長の情報に対して発声の仕方や癖による影響の補正が行われて、より話者に適合した特徴が抽出されるのである。その際に、話者の発声の仕方や癖

による影響が予め補正されている。したがって、より少量の発声データから良質の話者特徴が抽出される。

【0038】また、第7の発明のプログラム記録媒体は、コンピュータを、上記第1の発明の伸縮係数取得手段として機能させる話者特徴抽出処理プログラムが記録されていることを特徴としている。

【0039】上記構成によれば、上記第1の発明の場合と同様に、生理的な特徴の変動要因である声道長の情報に対して発声の仕方や癖による影響の補正が行われて、より話者に適合した特徴が抽出される。その際に、話者の発声の仕方や癖による影響が予め補正されているため、より少量の発声データから良質の話者特徴が抽出される。

【0040】

【発明の実施の形態】以下、この発明を図示の実施の形態により詳細に説明する。

<第1実施の形態>図1は、本実施の形態の音声認識装置におけるブロック図である。尚、この音声認識装置は、話者クラスタリング方式を用いた音声認識装置である。音声入力部1において、マイクから入力された音声はデジタル波形に変換されて音響分析部2および話者クラスタ選択部3に入力される。音響分析部2は、入力されたデジタル波形を短い時間間隔(フレーム)毎に周波数分析し、スペクトルを表す音響パラメータのベクトル系列に変換する。ここで、上記周波数分析には、MFCC(メル周波数F/T(高速フーリエ変換)ケプストラム)やLPC(線形予測分析)メルケプストラム等のスペクトルを効率よく表現できる方法が用いられる。こうして得られた音響パラメータ系列は、尤度(音韻類似度)演算部4に送出される。

【0041】上記話者クラスタ選択部3は、後に詳述する簡易音響モデル格納部6に話者クラスタ別に格納された簡易音響モデル(本実施の形態ではGMM(ガウシアン混合モデル)を使用)を用いて入力音声の話者クラスタを判定し、クラスタ情報を出力する。切換部5は、後に詳述する音響モデル格納部7に話者クラスタ別に格納された音響モデル(本実施の形態ではHMMを使用)の中から、話者クラスタ選択部3からのクラスタ情報に適合する話者クラスタの音響モデルを切り換え選択して尤度演算部4に送出する。そうすると、尤度演算部4は、音響分析部2からの入力音声の音響パラメータベクトルに対して切換部5からの音響モデルを作用させて、各音韻の状態毎に尤度を算出する。そして、得られた尤度系列を照合部8に送出する。

【0042】上記照合部8は、上記尤度演算部4からの尤度系列に対して、辞書格納部9に登録された総ての言語モデル(単語)との照合を行ない、各単語のスコアを算出する。そして、上位のスコアを呈する単語を認識候補(認識結果)として出力部10から出力するのである。

【0043】ここで、本実施の形態における音響モデル

格納部7は、学習話者のクラスタ数 n に応じて、第1音響モデル格納部7a、第2音響モデル格納部7b、…、第 n 音響モデル格納部7nの n 個の音響モデル格納部で構成されている。ここで、各音響モデル格納部7a~7nに格納される各音響モデルは、混合ガウス分布型のHMMである。この発明においては、生理的な特徴の変動に対処可能にすることを目的としており、話者性の大局的な安定要因である声道長をクラスタリング対象にするのである。尚、生理的な特徴の変動要因としては、上記声道長以外にも鼻腔、副鼻腔、声帯等の多くの要因があり、それらが絡み合って複雑な特徴を成している。したがって、個々の要因を数理的に扱うのは得策ではない。そこで、本実施の形態においては、各要因の複雑な特徴を混合ガウス分布型HMM(音響モデル)で表現するのである。

【0044】以下、上記音響モデル格納部7に格納される音響モデル群の作成方法について説明する。本実施の形態における音響モデル作成方法は、下記の(1)~(5)の5段階で構成される。

【0045】(1) 全話者モデルの作成

全学習話者の音声データを用いて学習を行い、混合数が1の音響モデル(全話者モデルと言う)を作成する。ここで、上記学習話者の集合を求めるに際して、男女別に2つの集合に分けてもよい。その場合には、大きく分けて男性用話者クラスタと女性用話者クラスタとの合計2種類の音響モデル群が生成されることになる。尚、上述のような全話者モデルに対して、後に第2実施の形態で述べる話者正規化方式による音声認識時において尤度演算部で用いる混合数が多数の音響モデルを、不特定話者モデルと呼ぶことにする。

【0046】(2) 特定話者モデルの作成

各学習話者毎に混合数が1の音響モデル(特定話者モデルと言う)を作成する。ここで、各学習話者に関しては、ある程度の量の音声データが整備されているものとする。ここで、1話者当りの音声データ量が多い場合には、HMMの学習アルゴリズムを用いて上記特定話者モデルを作成する。一方、1話者当りの音声データ量が少ない場合は、上記全話者モデルを基に上記VFSやMLLR(Maximum Likelihood Linear Regression)等の手法を用いて話者適応することによって上記特定話者モデルを作成する。

【0047】(3) 声道長正規化係数 α の導出

クラスタリングの基準となる声道長正規化係数 α を、補正係数 β と共に、各学習話者の特定話者モデルに関して、次に述べる非線形周波数ワーピング関数 $f()$ を用いて、式(1)に従って全話者モデルに対する尤度を最大にするという基準で求める。こうして、上記全話者モデルと特定話者モデルとの間の写像関係が、非線形周波数ワーピング関数 $f()$ を用いて求めることができる。

【0048】非線形周波数ワーピング関数 $f()$ ：

・ $x \leq \theta$ では、 $f(x) = \alpha \beta x$
 $(0.88 < \alpha < 1.13, 0.8 < \beta \leq 1)$ ($\theta = 1.5 \text{ kHz} \sim 1.8 \text{ kHz}$)

・ $\theta < x \leq \min(\omega/\alpha, \omega)$ では、 $f(x) = \alpha x$ ($\omega \approx 4 \text{ kHz}$)

・ $\min(\omega/\alpha, \omega) < x$ では、
 $\alpha > 1$ のとき $f(x) \rightarrow (\omega/\alpha, \omega)$ と $(fs/2, fs/2)$ とを結ぶ直線
 $\alpha \leq 1$ のとき $f(x) \rightarrow (\omega, \alpha \omega)$ と $(fs/2, fs/2)$ とを結ぶ直線

ここで、 θ : 広母音(日本語の場合「ア」や「オ」に相当)の第2フォルマントが存在する領域における上限周波数
 fs : サンプリング周波数

【0049】ここで、上記非線形周波数ワーピング関数 $f()$ における不連続性を考慮して、 θ 付近において直線「 $f(x) = \alpha \beta x$ 」と直線「 $f(x) = \alpha x$ 」とを接続させる。すなわち、 $(x, f(x))$ 座標上における $(k\theta, k\alpha\beta\theta)$ と $(\theta, \alpha\theta)$ とを直線で結ぶのである。 k は直線 $f(x) = \alpha \beta x$ 側の折れ線の頂点を表す定数であり、「0.7」程度の値とする。ここで、上記 θ は、男性の場合は 1.5 kHz 程度に設定され、女性の場合は 1.8 kHz 程度に設定される。尚、上記 θ は、話者の広母音における第2フォルマント周波数に応じて話者毎に設定を変えてもよい。また サンプリング周波数 fs は、本実施の形態においては 8 kHz 以上を仮定している。すなわち、 $fs = 12 \text{ kHz}$ の場合には、 $(fs/2, fs/2)$ は $(6 \text{ kHz}, 6 \text{ kHz})$ となるのである。

【0050】 $\alpha > 1$ である場合における上述のような折れ線で表される非線形周波数ワーピング関数 $f(x)$ を図2に示す。 α と β との変動範囲「 $0.88 < \alpha < 1.13$ 」, *30

$$\hat{\alpha} = \arg \max_{\alpha} \sum_{i \in \Omega} r_i(\mu_i^f) \quad (0.8 < \beta \leq 1) \quad \dots(1)$$

ここで、 Ω : 評価対象の出力確率密度関数集合のインデックス

$r_i()$: 標準モデルの第 i 番目の出力確率密度関数

$f()$: (α, β) を係数とする非線形周波数ワーピング関数

μ_i^f : 入力モデルにおける第 i 番目の出力確率密度関数

$b_i()$ の平均値ベクトル μ_i を $f()$ で周波数ワープしたベクトル。すなわち、 $C(f(C^{-1}(\mu_i)))$ となる。

C^{-1} , C : 音響パラメータからスペクトルへの変換とその逆変換

【0054】クラスタリングは、上記声道長正規化係数 α のみの1次元空間でのクラスタリングと、声道長正規化係数 α と補正係数 β との2次元空間でのクラスタリングとの2通りがある。声道長正規化係数 α は話者毎に固定されるが、補正係数 β は話者内で固定される場合と発話間で固定(=話者内で変動)される場合の2通りがある。どの範囲で補正係数 β を固定するかは話者の発声の癖に依存するため、話者に応じて使い分けるものとする。本実施の形態においては、音響モデルの各状態毎に

*「 $0.8 < \beta \leq 1$ 」は、夫々声道長の分布と声門上部の狭めにより生じるフォルマントの上昇度合いの観測結果に基づいて定めている。発声の仕方や癖によって声門上部の狭めが生じると、広母音の第2フォルマントの存在領域以下の低域のフォルマント周波数が上昇する。そのために、声道長正規化係数 α のみでは正しい声道長に写像できないことがある。上記係数 β は、この低域のフォルマント周波数の上昇に対する補正項なのである。

【0051】上記音響分析部2による音響分析で得られる音響パラメータや上記音響モデルの出力確率密度関数の引数は、通常MFCCやLPCケプストラムである。これらの音響パラメータの各次元はケプストラムと呼ばれる物理量であって、周波数ではない。そこで、上記周波数ワープ処理を行なう際には、学習データである音響パラメータからスペクトルへの変換 C^{-1} (ケプストラムの場合は逆cos変換)を行なって周波数次元に変換する。そして、周波数ワープ処理終了後は、逆変換 C (ケプストラムの場合はcos変換)を行なって元の音響パラメータ次元に戻すのである。

【0052】ここで、標準モデルと入力モデルとの2つの音響モデルの対応する状態間の尤度を、標準モデルの出力確率密度関数 $r_i()$ に、入力モデルの出力確率密度関数 $b_i()$ の平均値ベクトル μ_i を非線形周波数ワーピング関数 $f()$ で周波数ワープして得られたベクトルを代入したときの値と定義する。上記各出力確率密度関数は多次元ガウス分布であって、平均値ベクトルと分散ベクトルから成っている。

【0053】以上のことから、上記声道長正規化係数 α は式(1)によって最尤推定できることになる。

固定するものとする。

【0055】(4) 学習話者のクラスタリング

上記声道長正規化係数 α を用いて学習話者を所望の数 n にクラスタリングし、各学習話者の夫々が何れの話者クラスに属するかを決める。ここで、上記クラスタリングの方法は種々提案されているが、声道長正規化係数 α のみでクラスタリングする場合は、1次元空間であるから α 軸を n 個に等分割すればよい。尚、総学習話者数が少ない場合には、各話者クラスに属する話者数が等しくなるように n 個に分割してもよい。声道長正規化係数 α と補正係数 β との2次元空間上でクラスタリングする場合は、学習話者を k -means法等の手法によってクラスタリングすればよい。

【0056】ところで、何れの音響モデルの場合も、各話者クラスタの境界領域では学習サンプルが希薄だったり段差ができていたりしているために上手く学習されていない。したがって、入力話者が各話者クラスタの境界付近に位置する場合には、認識率が劣化するという(hard de

cision問題」が生じる。そこで、本実施の形態においては、この「hard decision問題」の対策として、学習話者を単純に分割するだけでなくオーバーラップさせて分割するのである。このオーバーラップは声道長に対応しているので物理的にも意味がある。すなわち、先ず n_0 個の話者クラスタに初期分割した後、 n_0 個の話者クラスタの各境界を中心として初期分割された話者クラスタにオーバーラップする($n_0 - 1$)個の話者クラスタに分割するのである。したがって、話者クラスタの数は合計($2n_0 - 1$)個となる。図3に、初期分割数 n_0 が「5」の場合のクラスタリング例を示す。縦軸は学習話者の頻度であり、横軸は声道長正規化係数 α である。5個の初期分割クラスタの境界を埋めるオーバーラップクラスタの数は4個であるから、総話者クラスタ数は合計9個となる。

【0057】上記オーバーラップさせるクラスタリングにおいて、オーバーラップのさせ方として、初期分割数 n_0 の異なる話者クラスタを併用してもよい。さらに、分割なしの全話者クラスタや男女別話者クラスタを併用してもよい。例えば、初期分割数 $n_0 = 7$ の話者クラスタに、初期分割数 $n_0 = 5$ の話者クラスタと男女別話者クラスタとを併用すると、 $(7 + 6) + (5 + 4) + 2$ の合計24個の話者クラスタとなる。

【0058】(5) 話者クラスタ別に音響モデルを作成上記 n 個の話者クラスタに属する総ての学習話者の音声データを用いて学習を行い、話者クラスタ毎に混合ガウス分布型HMMの音響モデルを作成する。話者クラスタへの初期分割数を n_0 個とすると、上記オーバーラップクラスタリングによって合計 $n = (2n_0 - 1)$ 個の音響モデルが生成されることになる。こうして作成された n 個の音響モデルの夫々が、音響モデル格納部7を構成する n 個の音響モデル格納部7a~7nの何れかに格納されるのである。尚、音響モデル格納部7a~7nの夫々に格納される「1個の音響モデル」とは、文字通り1つの音素の音響モデルを意味するのではなく、全音素に関する音響モデルの総称であることは言うまでも無い。

【0059】次に、上記簡易音響モデル格納部6について説明する。簡易音響モデル格納部6は、話者のクラスタ数 n に応じて、第1クラスタ用格納部6a、第2クラスタ用格納部6b、…、第 n クラスタ用格納部6nの n 個の簡易音響モデル格納部で構成されている。ここで、各クラスタ用格納部6a~6nに格納される各簡易音響モデルはGMMである。尚、GMMは、全音素を1状態で表す多混合連続分布型音響モデルである。

【0060】そして、上記話者クラスタ選択部3は音響分析手段を内蔵しており、入力音声から抽出された音響パラメータ系列に対して各クラスタ用格納部6a~6nに格納された総てのGMMを作用させて各GMM毎の尤度を算出する。そして、最も大きい尤度を呈するGMMが格納されたクラスタ用格納部6a~6nを表すクラスタ情

報を出力するのである。その場合、入力音声の正解音素列をユーザが教える必要がなく、教師なしで話者クラスタを選択することができる。すなわち、エンロールモードがないシステムにおいて有効なのである。

【0061】ここで、上記話者クラスタ選択の方法には、以下の[a]~[c]に示す3通りの方法がある。本実施の形態においては[b]の方法を用いている。

[a] 話者クラスタ音響モデル自身の利用

[b] 簡易型音響モデルの利用

10 [c] 声道長正規化係数 α および補正係数 β の直接推定

【0062】上記[a]の方法は、上記話者クラスタの音響モデル自身の尤度を用いる方法である。入力音声に対して教師語彙が与えられ、各話者クラスタにおける教師語彙の音響モデルを用いて認識処理を行い、各話者クラスタ毎の尤度を算出する。そして、最も大きい尤度を呈する話者クラスタを選択するのである。この選択方法は、エンロールモードにおいて入力音声の正解音素列をユーザが教えるという教師あり選択を基本としている。認識処理と同じ高精度な音響モデルを用いるので計算量は多くなるがエンロールによって正確なクラスタ選択が可能となる。

20 【0063】また、上記[c]の方法は、上述した音響モデルの作成方法における(3)の声道長正規化係数 α の導出で説明した手法と同様の手法を用いる。但し、特定話者音響モデルからではなく入力音声データから直接求めることになる。つまり、入力音声データに非線形周波数ワーピング関数 $f()$ を作用させて、全話者モデルを用いて最尤推定する方法で声道長正規化係数 α と補正係数 β とを求めるのである。この選択方法は、[a]や[b]の選択方法に比して不安定ではあるが、エンロールが可能であり、入力音声サンプルが多量にある場合には有効である。これは、後に第2実施の形態において説明する話者正規化で用いる手法と同じである。

【0064】上記構成において、入力音声の認識時には以下のように動作する。まず、話者クラスタ選択部3によって、上述のようにして最適な話者クラスタが選択され、クラスタ情報が切換部5に送出される。次に、尤度演算部4によって、切換部5で切り換え選択された話者クラスタの音響モデルを用いて尤度演算が行われ、得られた尤度系列が照合部8に送出される。そして、照合部8によって、ビタビサーチ等の探索アルゴリズムが用いられて辞書格納部9の言語モデルとの照合が行われ、各単語のスコアが算出される。尚、本実施の形態においては、照合部8による照合処理の前段処理が訴求点であるから、照合部8に関する詳細な説明は省略する。

【0065】上述のように、本実施の形態においては、上記音響モデル格納部7に格納する音響モデル群の作成に当って、学習話者を n 個の話者クラスタにクラスタリングする。そして、各話者クラスタに属する学習話者の音声データを用いた学習によって音響モデルを作成し、

各話者クラスタ別に第1音響モデル格納部7a〜第n音響モデル格納部7nに格納するようにしている。

【0066】その場合、上記学習話者のクラスタリングに際しては、各学習話者間の距離として、上記非線形周波数ワーピング関数 $f()$ を用いて、全学習話者の音響モデルに対する学習話者の音響モデルの尤度を最大にするという基準に従って求めた周波数軸の声道長正規化係数 α を用いるのである。さらに、発声の仕方や癖によって声門上部の狭めが生じると、広母音の第2フォルマントの存在領域以下の低域のフォルマント周波数が上昇する。そのために、声道長正規化係数 α のみでは正しい声道長に写像できないことがある。そこで、上記非線形周波数ワーピング関数 $f()$ に、上記低域のフォルマント周波数の上昇に対する補正項としての補正係数 β を導入している。

【0067】そして、標準モデル(全学習話者の音響モデル)と入力モデル(学習話者の音響モデル)の2つの音響モデルにおける対応する状態間の尤度を、標準モデルの出力確率密度関数 $r_i()$ に、入力モデルの出力確率密度関数 $b_i()$ の平均値ベクトル μ_i を非線形周波数ワーピング関数 $f()$ で周波数ワープして得られたベクトルを代入したときの値と定義して、上記声道長正規化係数 α を上記式(1)によって最尤推定するようにしている。

【0068】すなわち、本実施の形態によれば、上記学習話者のクラスタリング時に用いる各学習話者間の距離を、生理的な特徴の変動要因である声道長の情報と発声の仕方や癖による影響の補正情報とに基づいて設定することができる。したがって、より少量の発声データから発話者の癖を考慮した正確な各学習話者間の距離に基づいて、学習話者をクラスタリングできるのである。

【0069】また、実際の学習話者のクラスタリングに際しては、先ず n_0 個の話者クラスタに初期分割し、次に n_0 個の話者クラスタの各境界を中心として上記初期分割された話者クラスタにオーバーラップさせて $(n_0 - 1)$ 個の話者クラスタに分割し、合計 $n = (2n_0 - 1)$ 個の話者クラスタにクラスタリングするようにしている。したがって、各話者クラスタの境界領域では学習サンプルが希薄だったり段差ができていたりしているために上手く学習されず、認識率が劣化するという「harddecision問題」を解消することができるのである。

【0070】以上のことより、上述のようにしてクラスタリングされた各話者クラスタに属する学習話者の音声データ別に求められた音響モデルを上記第1音響モデル格納部7a〜第n音響モデル格納部7nに格納することによって、尤度演算部4は、より入力話者に適合した話者クラスタの音響モデルを適用することができる。したがって、高い認識率を得ることができるのである。

【0071】尚、上記実施の形態においては、上記話者クラスタ選択部3によって最適な話者クラスタを一つ選択するようにしているが、最適な話者クラスタを含む上

位複数の話者クラスタを選択するようにしてもよい。例えば、尤度の上位から k 個の話者クラスタを選択するとすると、そうすると、切換部5によって切り換え選択された k 個の音響モデルの夫々に関して、尤度演算部4によって尤度演算が行われて、照合部8に k 個の尤度系列が送られることになる。したがって、照合部8では、夫々の尤度系列に関して照合処理が行なわれ、最も大きい尤度を呈する単語/単語列が認識結果となるのである。

【0072】また、音声認識装置のハードウェア規模が大きく、計算量が許すのであれば、話者クラスタ選択部3による話者クラスタ選択を行わず、尤度演算部において総ての話者クラスタの音響モデルを用いて尤度演算処理を実行するようにしてもよい。この場合、各音響モデルを適用して得られた尤度が最大値を呈する単語/単語列が認識結果となる。

【0073】<第2実施の形態>図4は、本実施の形態の音声認識装置におけるブロック図である。尚、この音声認識装置は、話者正規化方式を用いた音声認識装置である。音声入力部11、音響分析部12、尤度演算部14、照合部18、辞書格納部19および出力部20は、図1に示す上記第1実施の形態における音声入力部1、音響分析部2、尤度演算部4、照合部8、辞書格納部9および出力部10と同様である。

【0074】周波数ワープ関数推定部15は、全話者音響モデル格納部16に格納された混合数が1の全話者モデル(HMM)を用いて、上記第1実施の形態における話者クラスタ選択方法[c]で述べたように、音響モデルの作成方法における(3)の声道長正規化係数 α の導出で説明した手法と同様の手法を用いて、入力音声データから非線形周波数ワーピング関数 $f()$ の声道長正規化係数 α および補正係数 β を推定する。そして、推定された声道長正規化係数 α および補正係数 β は、周波数ワープ部13に送出される。尚、全話者音響モデル格納部16に格納された全話者モデルは、上記第1実施の形態の音響モデルの作成方法における(1)の全話者モデルの作成で説明した手法と同様の手法で作成される。

【0075】そうすると、上記周波数ワープ部13は、上記推定値 (α, β) を係数とする非線形周波数ワーピング関数 $f()$ を用いて、入力音声の音響パラメータ系列を周波数ワープ(話者正規化)し、周波数ワープ後の音響パラメータ系列を尤度演算部14に送出するのである。そして、尤度演算部14では、周波数ワープされた音響パラメータ系列に対して、不特定話者音響モデル格納部17に格納された不特定話者モデル(HMM)を作用させて、各音韻の状態毎に尤度を算出するのである。

【0076】ところで、上記周波数ワープ関数推定部15における上記係数 (α, β) の推定方法には、以下に述べる二通りの推定方法がある。

(A) 入力音声データを直接用いる。

(B) 標準話者の音響モデルを入力音声データに話者適

応させた適応音響モデルを用いる。

そして、この二通りの推定方法を、入力音声データの量や質に応じて使い分けるのである。ここで、音声データの質とは尤度の上昇具合であり、周波数ワープ関数推定部 15 は、上記二通りの推定方法による尤度の上昇具合を見計らって、上昇の大きい推定方法を採用するのである。長いエンロール期間を許容できる音声認識装置の場合には、このような推定処理も可能となる。尚、長いエンロール期間を許容できない場合には、予め何れかの推定方法に固定しておけばよい。

【0077】上記推定方法(A)は、入力音声データが多い場合に有効であり、入力音声データから直接求めるために、精密な推定が可能となる。但し、入力音声データが少ない場合には、当該推定をエンロールモードで行う際に入力音声データに無い音素環境における係数(α , β)の推定や平滑化が問題になる。また、推定方法(B)は、入力音声データが少ない場合に有効であり、適応音響モデルの状態毎に補正係数 β を制御できるというメリットがある。

【0078】また、上記推定方法(A),(B)の各々に関して、使用する音響モデルは、全話者モデルの場合と、話者クラス別に作成された混合数が1の音響モデルの場合との二通りがある。音声認識装置の記憶容量が少ない場合には前者を採用する。一方、記憶容量が多い場合は音響モデル群を各話者クラス別に格納できるので後者を採用する。後者の場合には、入力音声データに基づいて最適な話者クラスを選択し、この選択話者クラスに属する音響モデルを使用することになる。すなわち、図4に示す音声認識装置は、全話者モデルを用いた推定方法(A)によって係数(α , β)の推定を行うのである。

【0079】以上、上記HMMに代表される音響モデルを用いた音声認識装置を例に、本実施の形態を説明したが、標準パターンとして音声波形または音響パラメータ系列を登録しておく音声認識装置に対しても、本実施の形態における話者正規化方法を適用することができる。その場合には、入力音声の音響パラメータ系列で成る特徴パターンと上記標準パターンとのマッチングには、上記HMMの場合の尤度に代ってスペクトル間の距離尺度を用いる。尚、その場合におけるマッチング部による処理手順を図5のフローチャートに示す。以下、図5に従って、標準パターンを登録しておく音声認識装置における上記マッチング部による処理手順について説明する。尚、この場合、係数(α , β)の更新幅と最大値とが予め設定されているものとする。

【0080】ステップS1で、上記特徴パターンと標準パターンの各フレーム間の対応関係(マッチングパスと言う)がDPマッチングによって求められる。その場合、上記DPマッチングに際しては、距離尺度としてケプストラム距離等のスペクトル間距離が用いられる。さ

らに、係数(α , β)に初期値が代入される。ステップS2で、上記マッチングパスにおける格子点位置が最初の格子点に初期化される。

【0081】ステップS3で、当該格子点に対応するフレームにおける入力音声の音響パラメータ系列に、(α , β)を係数とする非線形周波数ワーピング関数 $f()$ を作用させる。こうして、入力音声の当該フレームが係数(α , β)で非線形周波数ワープされる。

【0082】ステップS4で、当該格子点に対応するフレームにおける非線形周波数ワープ後の入力音声の特徴パターンと上記標準パターンとの累積距離が算出される。ステップS5で、次の格子点が在るか否かが判別される。その結果、在ればステップS6に進み、無ければステップS7に進む。ステップS6で、上記格子点位置が次の格子点に更新される。そうした後、上記ステップS3に戻って、次の格子点での処理に移行する。ステップS7で、上記係数(α , β)は上記最大値であるか否かが判別される。その結果、最大値であればステップS9に進む一方、そうでなければステップS8に進む。ステップS8で、係数(α , β)の値が上記更新幅だけ更新される。そうした後に、上記ステップS2に戻って、次の係数(α , β)での処理に移行する。このようにして、係数(α , β)の値を上記更新幅だけ順次更新しながら、上記非線形周波数ワープ後の入力音声と標準パターンとにおける上記マッチングパスに沿った累積距離が算出される。そして、上記ステップS7において、上記係数(α , β)は上記最大値であると判別されると上記ステップS9に進むのである。

【0083】ステップS9で、上記総ての累積距離の算出結果に基づいて、上記非線形周波数ワープ後の入力音声の特徴パターンと標準パターンとの累積距離を最小にする係数(α , β)の値が係数(α , β)の推定値として求められる。ステップS10で、上記係数(α , β)の推定値を係数とする上記非線形周波数ワーピング関数 $f()$ を用いて、入力音声の音響パラメータ系列を周波数ワープさせる。そして、周波数ワープ後の入力音声の音響パラメータ系列を照合部へ送出して、上記マッチング部による処理を終了するのである。

【0084】尚、上記マッチング部による処理の説明においては、単純なやり方で係数(α , β)の全部の組合せに関して累積距離を求めているが、山登り法や最急降下法等の高速に収束させる方法を採用しても差し支えない。

【0085】上述のように、本実施の形態においては、上記周波数ワープ関数推定部15によって、入力音声の音響パラメータ系列から非線形周波数ワーピング関数 $f()$ の係数(α , β)を推定する。そして、周波数ワープ部13によって、上記推定値(α , β)を係数とする非線形周波数ワーピング関数 $f()$ を用いて入力音声の音響パラメータ系列を周波数ワープすることによって、話者正規

化するようにしている。

【0086】その場合、発声の仕方や癖によって声門上部の狭めが生じると、広母音の第2フォルマントの存在領域以下の低域のフォルマント周波数が上昇する。そのために、上記声道長正規化係数 α のみを用いた非線形周波数ワーピング関数 $f()$ では正しい声道長に写像できないことがある。そこで、上記非線形周波数ワーピング関数 $f()$ に、上記低域のフォルマント周波数の上昇に対する補正項としての補正係数 β を導入している。

【0087】そして、全学習話者の音響モデルと入力音声の音響モデルとの2つの音響モデルにおける対応する状態間の尤度を、標準モデルの出力確率密度関数 $r_i()$ に、入力モデルの出力確率密度関数 $b_i()$ の平均値ベクトル μ_i を非線形周波数ワーピング関数 $f()$ で周波数ワーピングして得られたベクトルを代入したときの値と定義して、上記係数 (α, β) を上記式(1)によって最尤推定するようにしている。

【0088】すなわち、本実施の形態によれば、上記話者正規化する際に用いる非線形周波数ワーピング関数 $f()$ の係数 (α, β) として、生理的な特徴の変動要因である声道長の情報である声道長正規化係数 α に対して発声の仕方や癖による影響の補正を行ったものを用いることができる。したがって、発話者の癖を考慮した上記非線形周波数ワーピング関数 $f()$ に基づいて、より少量の発声データから、標準話者のスペクトルにより近い周波数特性を有するように入力音声を話者正規化できるのである。

【0089】<第3実施の形態>図6は、本実施の形態の音声認識装置におけるブロック図である。尚、この音声認識装置は、話者適応方式を用いた音声認識装置である。音声入力部21、音響分析部22、尤度演算部27、照合部29、辞書格納部30及び出力部31は、図1に示す上記第1実施の形態における音声入力部1、音響分析部2、尤度演算部4、照合部8、辞書格納部9および出力部10と同様である。また、周波数ワーピング関数推定部23、全話者音響モデル格納部24および不特定話者音響モデル格納部26は、図4に示す上記第2実施の形態における周波数ワーピング関数推定部15、全話者音響モデル格納部16および不特定話者音響モデル格納部17と同様である。

【0090】すなわち、上記周波数ワーピング関数推定部23は、上記第2実施の形態の場合と同様にして、全話者モデルを用いて、入力音声データから非線形周波数ワーピング関数 $f()$ の係数 (α, β) を推定する。そして、周波数ワーピング部25によって、この推定された係数 (α, β) の逆数を係数とする非線形周波数ワーピング関数 $f()$ を用いて、不特定話者音響モデル格納部26に格納された不特定話者モデルを周波数ワーピングする。こうして周波数ワーピングされた不特定話者音響モデルを、話者適応モデル(HMM)として話者適応音響モデル格納部28に格納す

る。そうすると、尤度演算部27は、音響分析部22からの入力音声の音響パラメータ系列に対して、話者適応音響モデル格納部28に格納された話者適応モデルを作用させて、上述した尤度演算処理を行なうのである。

【0091】このように、本実施の形態においては、上記周波数ワーピング関数推定部23によって、入力音声の音響パラメータ系列から非線形周波数ワーピング関数 $f()$ の係数 (α, β) を推定する。そして、周波数ワーピング部25によって、上記推定値 (α, β) の逆数を係数とする非線形周波数ワーピング関数 $f()$ を用いて不特定話者モデルを周波数ワーピングすることによって、不特定話者モデルを話者適応させるようにしている。

【0092】その場合、発声の仕方や癖によって声門上部の狭めが生じると、広母音の第2フォルマントの存在領域以下の低域のフォルマント周波数が上昇する。そのために、上記声道長正規化係数 α のみを用いた非線形周波数ワーピング関数 $f()$ では正しい声道長に写像できないことがある。そこで、上記非線形周波数ワーピング関数 $f()$ に、上記低域のフォルマント周波数の上昇に対する補正項としての補正係数 β を導入している。

【0093】そして、全学習話者の音響モデルと入力音声の音響モデルとの2つの音響モデルにおける対応する状態間の尤度を、標準モデルの出力確率密度関数 $r_i()$ に、入力モデルの出力確率密度関数 $b_i()$ の平均値ベクトル μ_i を非線形周波数ワーピング関数 $f()$ で周波数ワーピングして得られたベクトルを代入したときの値と定義して、上記係数 (α, β) を上記式(1)によって最尤推定するようにしている。

【0094】すなわち、本実施の形態によれば、上記不特定話者モデルを話者適応する際に用いる非線形周波数ワーピング関数 $f()$ の係数 (α, β) として、生理的な特徴の変動要因である声道長の情報である声道長正規化係数 α に対して発声の仕方や癖による影響の補正を行ったものを用いることができる。したがって、発話者の癖を考慮した上記非線形周波数ワーピング関数 $f()$ に基づいて、より少量の発声データから入力話者のスペクトルにより近い周波数特性を有するように不特定話者モデルを話者適応できるのである。

【0095】尚、本実施の形態における上記話者適応音響モデル格納部28に格納する話者適応モデルの与え方には、上述の与え方の以外に、話者クラスタを用いる方法を採用してもよい。そして、この二通りの与え方を、音声認識装置の規模や入力音声データの量や質に応じて使い分けるのである。ここで、音声データの質とは尤度の上昇具合であり、周波数ワーピング関数推定部23は、上記二通りの与え方による尤度の上昇具合を見計らって、上昇の大きい推定方法を採用するのである。長いエンロール期間が許容できる音声認識装置の場合には、このような推定処理も可能となる。尚、上記話者クラスタを用

形態における選択法[a]に基づいて話者クラスタを選択する。そして、選択された話者クラスタの音響モデルを話者適応モデルとして話者適応音響モデル格納部28に格納するのである。

【0096】また、上述した二つの与え方の何れかによって得られた話者適応モデルを初期モデルとして、上記MLLR方やVFS法等の既存の話者適応技術を用いて話者適応を行って新たに話者適応モデルを生成し、これを尤度演算部で用いるようにしても差し支えない。尚、その際における話者クラスタのクラスタ数や補正係数 β の使い方(固定範囲)やクラスタ選択の方法等は、適用する音声認識装置あるいは用いる音響モデルの規模やエンロールモードの有無等に依存するため、それらの項目については、本実施の形態においては規定しない。例えば、コンパクトな音声認識装置を望む場合には、上記話者クラスタの数は減ることになる。また、音響モデルの規模が小さい場合には、補正係数 β は状態毎に固定すればよい。エンロールモードが許容できない場合には、GMMを用いたクラスタ選択が行われることになる。

【0097】<第4実施の形態>図7は、本実施の形態のテキスト音声合成装置におけるブロック図である。なお、このテキスト音声合成装置は、声質変換方式を用いたテキスト音声合成装置である。テキスト解析部41は、単語とそのアクセント型とが格納されたアクセント辞書42を用い、入力テキストに対して形態素解析および係り受け解析を行って音素文字列とアクセント情報とを生成して韻律生成部43に送出する。韻律生成部43は、韻律制御テーブル44を参照して、継続時間長やピッチやパワーの韻律情報を生成して、音素文字列と共に音声素片選択部45に送出する。そうすると、音声素片選択部45は、音声素片辞書46から音素環境や韻律環境に最適な音声素片を選択し、音声素片情報を生成する。そして、この生成された音声素片情報を周波数ワープ部48に出力する一方、上記韻律情報を音声素片合成部47に出力する。

【0098】一方、周波数ワープ関数推定部49は、声質変換のターゲット話者の入力音声波形を基に、第2、第3実施の形態の場合と同様にして、上記非線形周波数ワッピング関数 $f()$ の係数 (α, β) を推定する。そうすると、周波数ワープ部48は、この推定された係数 (α, β) の逆数を係数とする非線形周波数ワッピング関数 $f()$ を用いて上記音声素片情報である音響パラメータ系列を周波数ワープし、周波数ワープ後の音声素片情報を音声素片合成部47に送出する。最後に、音声素片合成部47は、周波数ワープ部48からの周波数ワープ後の音声素片情報(音声素片の音響パラメータ系列)と音声素片選択部45からの韻律情報とを用いて、音声波形を生成しスピーカ50から音声出力する。

【0099】上述のように、本実施の形態においては、テキスト音声合成を行うに際して、上記周波数ワープ関

数推定部49によって、声質変換のターゲット話者における入力音声の音響パラメータ系列から非線形周波数ワッピング関数 $f()$ の係数 (α, β) を推定する。そして、周波数ワープ部48によって、上記推定値 (α, β) を係数とする非線形周波数ワッピング関数 $f()$ を用いて、テキストに基づいて選択された音声素片の音響パラメータ系列を周波数ワープすることによって、声質変換を行うようにしている。

【0100】その場合、上記係数 (α, β) を推定に際しては、発声の仕方や癖によって声門上部の狭めが生じると、広母音の第2フォルマントの存在領域以下の低域のフォルマント周波数が上昇する。そのために、上記非線形周波数ワッピング関数 $f()$ に、上記低域のフォルマント周波数の上昇に対する補正項としての補正係数 β を導入している。

【0101】そして、全学習話者の音響モデルと入力音声の音響モデルとの2つの音響モデルにおける対応する状態間の尤度を、標準モデルの出力確率密度関数 $r_i()$ に、入力モデルの出力確率密度関数 $b_i()$ の平均値ベクトル μ_i を非線形周波数ワッピング関数 $f()$ で周波数ワープして得られたベクトルを代入したときの値と定義して、上記声道長正規化係数 α を、上記式(1)によって最尤推定するようにしている。

【0102】すなわち、本実施の形態によれば、上記声質変換を行う際に用いる非線形周波数ワッピング関数 $f()$ の係数 (α, β) として、生理的な特徴の変動要因である声道長の情報である声道長正規化係数 α に対して発声の仕方や癖による影響の補正を行ったものを用いることができる。したがって、発話者の癖を考慮した上記非線形周波数ワッピング関数 $f()$ に基づいて、より少量の発声データから、ターゲット話者のスペクトルにより近い周波数特性を有するように音声素片情報を声質変換できるのである。

【0103】本実施の形態はスペクトル包絡の変換であり、声質の適応におおいに効果がある。しかしながら、話者間の声の特徴差は声質だけではなく韻律が大きく寄与する。したがって、本実施の形態に対して韻律の適応技術を併用しても構わない。

【0104】尚、上述した各実施の形態においては、上記声道長正規化係数 α と補正係数 β とで成る話者特徴を用いてクラスタリングされた音響モデルを搭載した音声認識装置、上記声道長正規化係数 α と補正係数 β とで成る話者特徴を用いて話者正規化あるいは話者適応を行う音声認識装置、および、上記声道長正規化係数 α と補正係数 β とで成る話者特徴を用いて声質変換を行う音声合成装置について説明している。しかしながら、この発明は、上記声道長正規化係数 α と補正係数 β とを話者特徴として抽出する話者特徴抽出装置にも適用されるものである。

【0105】ところで、その場合の話者特徴抽出装置に

における上記伸縮係数取得手段としての機能は、プログラム記録媒体に記録された話者特徴抽出処理プログラムによって実現される。上記プログラム記録媒体は、ROM(リード・オンリ・メモリ)でなるプログラムメディアである。あるいは、外部補助記憶装置に装着されて読み出されるプログラムメディアであってもよい。尚、何れの場合においても、上記プログラムメディアから話者特徴抽出処理プログラムを読み出すプログラム読み出し手段は、上記プログラムメディアに直接アクセスして読み出す構成を有していてもよいし、RAM(ランダム・アクセス・メモリ)に設けられたプログラム記憶エリア(図示せず)にダウンロードして、上記プログラム記憶エリアにアクセスして読み出す構成を有していてもよい。尚、上記プログラムメディアからRAMの上記プログラム記憶エリアにダウンロードするためのダウンロードプログラムは、予め本体装置に格納されているものとする。

【0106】ここで、上記プログラムメディアとは、本体側と分離可能に構成され、磁気テープやカセットテープ等のテープ系、フロッピー(登録商標)ディスク、ハードディスク等の磁気ディスクやCD(コンパクトディスク)-ROM, MO(光磁気)ディスク, MD(ミニディスク), DVD(デジタルビデオディスク)等の光ディスクのディスク系、IC(集積回路)カードや光カード等のカード系、マスクROM, EPROM(紫外線消去型ROM), EEPROM(電気的消去型ROM), フラッシュROM等の半導体メモリ系を含めた、固定的にプログラムを担持する媒体である。

【0107】また、上記各実施の形態における音声認識装置、音声合成装置および話者特徴抽出装置は、モデムを備えてインターネットを含む通信ネットワークと接続可能な構成を有していれば、上記プログラムメディアは、通信ネットワークからのダウンロード等によって流動的にプログラムを担持する媒体であっても差し支えない。尚、その場合における上記通信ネットワークからダウンロードするためのダウンロードプログラムは、予め本体装置に格納されているものとする。または、別の記録媒体からインストールされるものとする。

【0108】尚、上記記録媒体に記録されるものはプログラムのみに限定されるものではなく、データも記録することが可能である。

【0109】

【発明の効果】以上より明らかなように、第1の発明の話者特徴抽出装置は、伸縮係数取得手段によって、音声のスペクトルにおける広母音の第2フォルマントの存在領域以下の低い周波数領域で伸縮係数 α に補正係数 β を乗じて補正を行った非線形周波数ワーピング関数を用いて、標準話者の音声パターンに対して入力話者の音声パターンの尤度を最大にするという基準に従って上記伸縮係数 α を求め、求めた伸縮係数 α を話者特徴とするので、生理的な特徴の変動要因である声道長の情報に対し

て発声の仕方や癖による影響の補正を行って、より話者に適合した特徴を抽出することができる。さらに、上記発声の仕方や癖を表す発声データを必要とはせず、より少量の発声データから精度良く話者特徴を抽出できる。

【0110】また、上記第1の発明の話者特徴抽出装置は、上記伸縮係数取得手段を、上記非線形周波数ワーピング関数を用いて、標準話者の音声パターンに対して入力話者の音声パターンの尤度を最大にするという基準に従って上記補正係数 β をも求めるように成せば、さらに話者に適合した特徴を抽出することができる。

【0111】また、第2の発明の音声認識装置は、音響モデルを、上記伸縮係数 α に基づいて学習話者をクラスタリングして得られた各話者クラスタ別に格納する際に、上記伸縮係数 α を、上記非線形周波数ワーピング関数を用いた最尤推定によって求めるので、生理的な特徴の変動要因である声道長の情報に対して発声の仕方や癖による影響の補正を行って、より学習話者の音声パターンに適合した学習話者間の距離を用いて上記クラスタリングを行うことができる。したがって、この発明によれば高い認識率を得ることができる。さらに、上記発声の仕方や癖を表す発声データを必要とはせず、より少量の発声データから学習話者間の距離を得ることができるのである。

【0112】また、上記第2の発明の音声認識装置は、上記学習話者のクラスタリングを上記伸縮係数 α と補正係数 β との2次元平面に対して行うようになっており、上記補正係数 β を上記非線形周波数ワーピング関数を用いた最尤推定によって求めれば、さらに学習話者の音声パターンに適合した距離を用いてクラスタリングを行うことができる。

【0113】また、上記第2の発明の音声認識装置は、上記話者クラスタを、所定のクラスタ数の初期クラスタと、上記各初期クラスタの境界を含んで上記各初期クラスタにオーバーラップするオーバーラップクラスタとで構成すれば、学習サンプルが希薄だったり段差がでやすい上記各初期クラスタの境界領域を、何れかのオーバーラップクラスタに属させることができる。したがって、上記各初期クラスタの境界領域において認識率が劣化するという「hard decision問題」を解消できる。

【0114】また、第3の発明の音声認識装置は、正規化手段を、上記非線形周波数ワーピング関数を用いて伸縮係数 α と補正係数 β とを最尤推定する周波数ワーピング関数推定手段と、上記推定された α と β とを係数とする上記非線形周波数ワーピング関数を用いて上記入力話者の音声スペクトルの周波数軸を伸縮する周波数ワープ手段で構成したので、より標準話者の音声スペクトルに近づくように話者を正規化することができる。したがって、この発明によれば高い認識率を得ることができる。さらに、上記発声の仕方や癖を表す発声データを必要とはせず、より少量の発声データに基づいて話者正規化を

行うことができる。

【0115】また、第4の発明の音声認識装置は、話者適応手段を、上記非線形周波数ワーピング関数を用いて伸縮係数 α と補正係数 β とを最尤推定する周波数ワーピング関数推定手段と、上記推定された α の逆数と β の逆数とを係数とする上記非線形周波数ワーピング関数を用いて音響モデルの周波数軸を伸縮する周波数ワープ手段で構成したので、より入力話者の音声スペクトルに近くように話者適応を行うことができる。したがって、この発明によれば、高い認識率を得ることができる。さらに、上記発声の仕方や癖を表す発声データを必要とせず、より少量の発声データに基づいて話者適応を行うことができる。

【0116】また、上記第3の発明あるいは第4の発明の音声認識装置は、上記周波数ワーピング関数推定手段を、上記入力話者の音声パターンの代わりに、標準話者の音響モデルを上記入力話者の音声パターンに話者適応させた適応音響モデルを用いるように成せば、入力話者の音声パターン数が少ない場合でも対処することができる。さらに、上記適応音響モデルの状態毎に補正係数 β を制御して、話者の発声の仕方や癖による入力音声パターンのずれを木目細かく補正することが可能になる。

【0117】また、上記第2の発明乃至第4の発明の何れか一つの発明の音声認識装置は、上記補正係数 β をサブワード単位に求め、上記サブワード毎に決定すれば、上記補正係数 β を上記サブワード単位で変更することができ、話者の発声の仕方や癖による入力音声パターンのずれを木目細かく補正することができる。

【0118】また、第5の発明の音声合成装置は、声質変換手段を、上記非線形周波数ワーピング関数を用いて伸縮係数 α と補正係数 β とを最尤推定する周波数ワーピング関数推定手段と、上記推定された α の逆数と β の逆数とを係数とする上記非線形周波数ワーピング関数を用いて標準話者の音声素片の周波数軸を伸縮する周波数ワープ手段で構成したので、より発話者の声質に適合するように合成音声の声質を変換することができる。さらに、上記発声の仕方や癖を表す発声データを必要とせず、より少量の発声データに基づいて声質変換を行うことができる。

【0119】また、上記第5の発明の音声合成装置は、上記周波数ワーピング関数推定手段をサブワード単位に上記補正係数 β を求め、上記サブワード毎に上記補正係数 β を推定するように成せば、上記補正係数 β を上記サブワード単位に変更することができ、発話者の発声の仕方や癖による入力音声パターンのずれを木目細かく補正することができる。

【0120】また、第6の発明の話者特徴抽出方法は、上記非線形周波数ワーピング関数を用いて、最尤度推定によって上記伸縮係数 α を求めて話者特徴とするので、

生理的な特徴の変動要因である声道長の情報に対して発声の仕方や癖による影響の補正を行って、より話者に適合した特徴を抽出することができる。さらに、上記発声の仕方や癖を表す発声データを必要とせず、より少量の発声データから精度良く話者特徴を抽出できる。

【0121】また、第7の発明のプログラム記録媒体は、コンピュータを、上記第1の発明における上記伸縮係数取得手段として機能させる話者特徴抽出処理プログラムが記録されているので、上記第1の発明の場合と同様に、より話者に適合した特徴を抽出することができる。さらに、より少量の発声データから良質の話者特徴を抽出できる。

【図面の簡単な説明】

【図1】 この発明の話者クラスタリング方式を用いた音声認識装置におけるブロック図である。

【図2】 非線形周波数ワーピング関数の一例を示す図である。

【図3】 初期分割数が5である場合のクラスタリング例を示す図である。

【図4】 図1とは異なる話者正規化方式を用いた音声認識装置のブロック図である。

【図5】 標準パターンを用いる音声認識装置に図4と同様の話者正規化方式を適用した際におけるマッチング部による処理手順のフローチャートである。

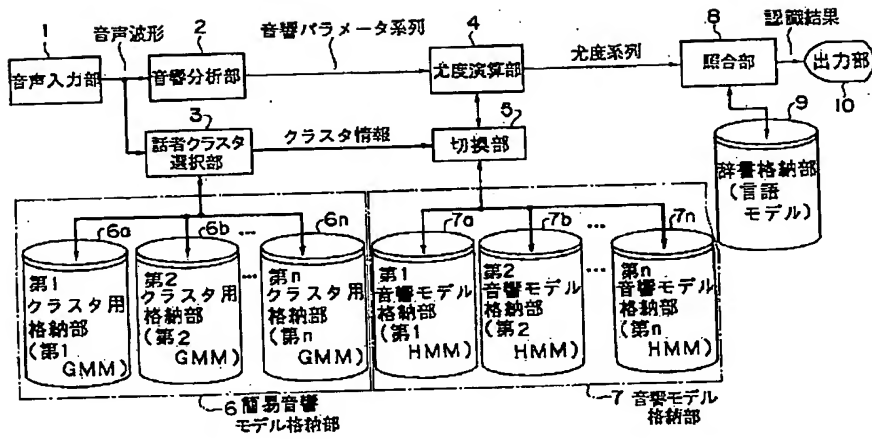
【図6】 図1および図4とは異なる話者適応方式を用いた音声認識装置におけるブロック図である。

【図7】 この発明の声質変換方式を用いた音声合成装置におけるブロック図である。

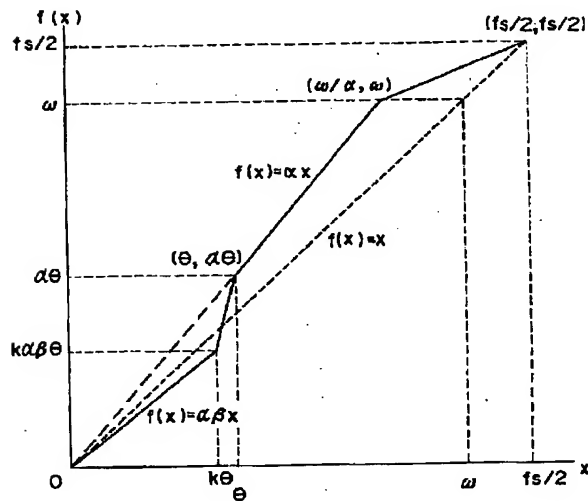
【符号の説明】

- 1, 11, 21…音声入力部、
- 2, 12, 22…音響分析部、
- 3…話者クラスタ選択部、
- 4, 14, 27…尤度演算部、
- 5…切換部、
- 6…簡易音響モデル格納部、
- 7…音響モデル格納部、
- 8, 18, 29…照合部、
- 9, 19, 30…辞書格納部、
- 10, 20, 31…出力部、
- 13, 25, 48…周波数ワープ部、
- 15, 23, 49…周波数ワープ関数推定部、
- 16, 24…全話者音響モデル格納部、
- 17, 26…不特定話者音響モデル格納部、
- 28…話者適応音響モデル格納部、
- 41…テキスト解析部、
- 43…韻律生成部、
- 45…音声素片選択部、
- 47…音声素片合成部、
- 50…スピーカ。

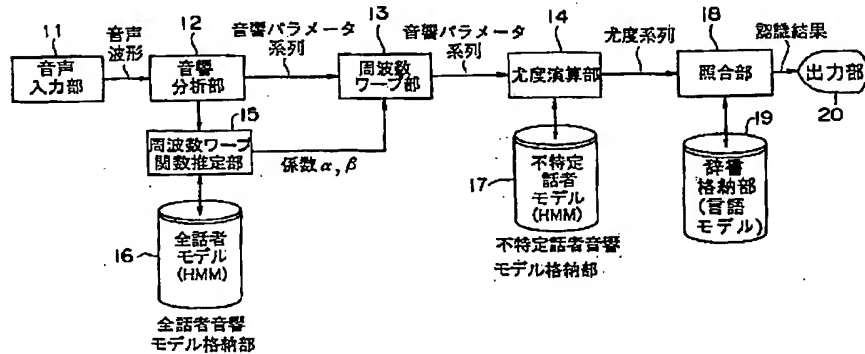
【図1】



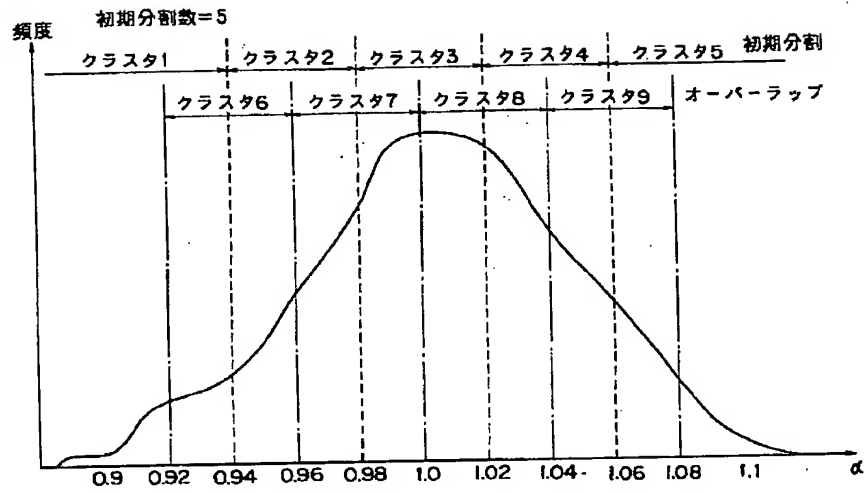
【図2】



【図4】



【図3】



【図5】

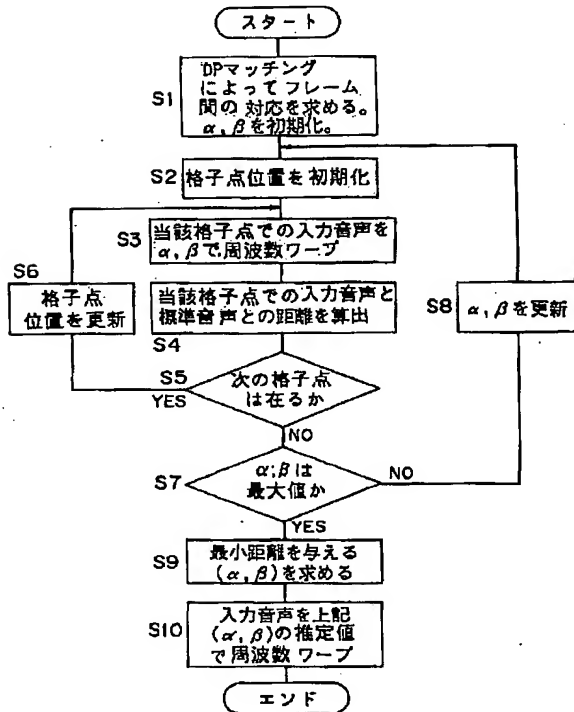


Figure 1 is a block diagram of a Japanese speech synthesis system. The system consists of the following components and data flows:

- Input Voice Waveform (49)**: Provides input to the **Pitch Contour Waveform Analysis/Estimation Unit (43)**.
- Input Text (41)**: Provides input to the **Text Analysis Unit (41)**.
- Pitch Contour Waveform Analysis/Estimation Unit (43)**: Extracts **Pitch Contour Information (45)** from the input waveform.
- Text Analysis Unit (41)**: Extracts **Phoneme and Pitch Information (43)** from the input text.
- Phoneme Selection Unit (45)**: Receives **Phoneme and Pitch Information (43)** and **Pitch Contour Information (45)** to select **Phoneme Elements (48)**.
- Phoneme Element Synthesis Unit (47)**: Receives **Phoneme Elements (48)** and **Pitch Contour Information (45)** to synthesize **Phoneme Elements (47)**.
- Phoneme Element Information (48)**: The output of the **Phoneme Selection Unit (45)**.
- Phoneme Elements (47)**: The output of the **Phoneme Element Synthesis Unit (47)**.
- Speech Waveform (50)**: The final output of the system, generated by a **Speaker (50)**.

301A